# The ESIS Query Environment Pilot Project N94-23952

Jens J. Fuchs[1]   Alessandro Ciarlo[2]   Stefano Benso[3]

[1]Computer Resources International A/S
Bregneroedvej 144, DK-3460 Denmark
Phone: +45 45822100

[2]ESA/ESRIN ISD
Via Galileo Galilei
I-00044 Frascati, Italy
Phone: +39 6 941801

[3]CISET Space Division
Via Salaria 1027
I-00138 Roma, Italy
Phone +39 6 881701

## 1. Introduction

The European Space Information System (ESIS) was originally conceived to provide the European space science community with simple and efficient access to space data archives, facilities with which to examine and analyze the retrieved data, and general information services. To achieve that ESIS will provide the scientists with a discipline specific environment for querying in a uniform and transparent manner data stored in geographically dispersed archives. Furthermore it it will provide discipline specific tools for displaying and analyzing the retrieved data.

The central concept of ESIS is to achieve a more efficient and wider usage of space scientific data, while maintaining the physical archives at the institutions which created them, and has the best background for ensuring and maintaining the scientific validity and interest of the data. In addition to coping with the physical distribution of data, ESIS is to manage also the heterogenity of the individual archives' data models, formats and data base management systems. Thus the ESIS system shall appear to the user as a single database, while it does in fact consist of a collection of dispersed and locally managed databases and data archives.

The work reported in this paper is one of the results of the ESIS Pilot Project which is to be completed in 1993. More specifically it presents the pilot ESIS Query Environment (ESIS QE) system which forms the data retrieval and data dissemination axis of the ESIS system. The others are formed by the ESIS Correlation Environment (ESIS CE) and the ESIS Information Services.

The ESIS QE Pilot Project is carried out for the European Space Agency's Research and Information center, ESRIN, by a Consortium consisting of Computer Resources International, Denmark, CISET S.p.a, Italy, the University of Strasbourg, France and the Rutherford Appleton Laboratories in the U.K. Furthermore numerous scientists both within ESA and space science community in Europe have been involved in defining the core concepts of the ESIS system.

## 2. The ESIS Query Environment

Past and future satellite missions have and will generate massive amounts of data which must be archived and indexed for subsequent scientific processing and analysis. Together with bibliographic databases and object catalogues these data archives represents an enormous and very valuable source of scientific information. However, the efficient use of these data banks are today obstructed by a number of factors:

o   The structures and tools available for querying and working with a given database are often highly mission-dependent, and not standardized.

o   Database management systems and the query languages provided are different making it difficult to learn how to use new database systems.

o   Data in different databases cannot easily be combined.

C-10

o   Instruments may have different sensitivity, filters, resolution etc.

o   A number of different data exchange formats exist for transfer of data to users equipment.

As the cost of a satellite mission is very high, it is highly desirable to be able to make efficient use of data already collected thereby enabling future satellite missions to have more specific goals, supplementing data already collected, or providing necessary additional information on certain objects.

The initial goal of the ESIS QE is to overcome the mentioned obstacles, by providing researchers with an efficient tool for access to and use of space scientific data. Furthermore the project shall provide future missions with a structured environment for making their data available to researchers. This will include standard formats, structures and query languages.

## 2.1 Interoperability of the ESIS QE

Providing common access to two or more databases which use different DBMS can be achieved either by trying to get them to work transparently together or by providing higher or lower degrees of interoperability.

The lowest level of interoperability may be achieved by a basic system containing a directory in which the user can identify the archives containing specific information, and a communication infrastructure allowing remote acces to the archives. The next levels are characterized by improved communication facilities providing automatic logon to remote hosts, and by the introduction of common query languages which can be used with the archive specific datamodels, and are interpreted by front end servers of the specific archives. The highest levels of interoperability introduces a common data model or schema, but may explicitly fail to address all the problems of syntactic and semantic consistency which must in principle be solved before the complete integration is achieved.

Within the ESIS QE Pilot Project the original aim was to achieve a very high level of interoperability through the provision of a common schema within each releavnt scientific discipline, and a mapping mechanism allowing these common schema to be mapped phyiscally and semantically on the underlying archives.

## 2.2   The Scientific Disciplines

In the pilot project two scientific domains were selected, namely:

o   Astronomy and Astrophysics

o   Space Physics

Within astronomy the following archives were selected for integration into the ESIS QE system:

o   IUE: (Infrared Ultraviolet Explorer, Villafranca Tracking Station, Spain)

o   STARCAT: (ESO, Garching, Germany)

o   EXOSAT: (ESTEC, Noordwijk, The Netherlands)

o   SIMBAD: (Strasbourg Observatory, France)

For Space Physics, only one archive was selected, namely:

o   GDF: (Geophysical Data Facility, RAL, United Kingdom).

Finally both domains were to be supported in the area of bibliographic data by providing access to the ESA-IRS database system operated at ESRIN.

It is foreseen that users will have varying degrees of computer literacy and database experience. It is therefore considered extremely important that the ESIS QE can cater to the occasional user, requiring help and assistance, as well as to the experienced user who may have his/her favorite database environment.

## 2.3   The QE Functions

The ESIS QE system provides the following main functionalities:

o   A user interface supporting the formulation of queries in various ways.

o   For each domain a common conceptual data model, defining the entities and attributes available for querying.

o   The translation of the query expressed relative to the conceptual model into a set of queries targeted against one or more archives, i.e. physical data-

models.

o The retrieval and integration of the physical data from the various archives, and the forwarding and presentation of these data back to the user.

o The provision of a user specific data repository, the Personal Database, allowing the user to maintain a collection of retrieved data for subsequent correlation and refinement.

o Facilities for import and export of data allowing the retrieved data to be analyzed further using tools of the individual scientist or e.g. the ESIS CE.

## 3. The ESIS QE Architecture

The ESIS QE system is a distributed system by nature. It must provide a geographically dispersed user community with access to an almost equally dispersed set of archives and databases in the most effective and cost efficient manner. To achieve this a layered client/server approach has been adopted such that conceptually the system can be viewed as being composed of three logical components and the physical Archive Layer, which work together transparently.

### Query Agent component (QA)

Provides the user interface and query formulation and result presentation facilities based on the common conceptual model. Provides service functions and interfaces to other systems, e.g. the Correlation Environment.

### Query Executer component

Collects the user queries identifies the potential archives which may be sources of data for the specified query. Generates a set of archive specific queries and forwards these to the next layer. Upstream it collects results from the archives, integrates them and forwards them to the user.

### Local Query Server component

Translates the archive specific query into the local archive query language and submits it for execution. Retrieves the results and homogenizes it depending on data type and use. Applies special functions not directly available in the archives on the temporary results and forwards the result to the Query Executer component.

Host Archive

The host database system which contains data such as catalogues of astronomical objects, mission logs, data sets, bibliographic references, general information about the scientific community.

The software components constituting the three top layers will be distributed among three physical platforms: the User Platform, the Access Point and the Archive Node. The User Platform hosts software responsible for the Query Agent functionality, the Access Point the Query Executer and the Archive Node the Local Query Server functionality. However it is foreseen that the same physical machine will host modules overlapping the conceptual components.

### 3.1 The Pilot System Architecture

Figure 1 provides an overview of the overall pilot system architecture. The hardware architecture is composed of several distributed computers interconnected via local area networks and wide are telecommunication networks. The constituents of this architecture are:

o The User Platform, which are the work stations used by the scientific user to access the ESIS system, and do local data processing. A User Platform will be located at the scientists institution, interconnected via a network to an Access Point. The ESIS QE pilot system will support three different types of workstations: VAX/VMS with DECWindows Motif, SUN/UNIX with OSF/Motif.

o The Access Point, is a VMS based machine which provides the first layer of server capability, and the further access to the Archive Nodes. In the pilot configuration there will be only one Access Point, located at ESRIN. However the overall concept is to have distributed Access Points, where each will serve a subset of users, in order to remove a potential bottleneck and reduce the communications overhead and cost incurred by one central node.

o The Archive Nodes, are front-end or interface computers of the data archives, hosting the LQS functions.

As indicated by the figure the pilot configuration provides access to the IUE, EXOSAT, HST, GDF and Simbad archives as was the original aim. The

ESA-IRS database have been replaced by the ESISBIB database and furthermore the ESIS Cats and Logs database has been added. Both these databases are located at ESRIN but may conceptually be seen as individual archives, each with their own dedicated server (Local Query Server).

The ESIS Cats & Logs (ESISCAT) database contains a partially homogenized collection of the astronomical catalogues and observation logs contained also in IUE, EXOSAT, and STARCAT/HST, plus a number of new databases developed within the ESIS project. ESISCAT was introduced into the originally fully distributed architecture in order to overcome various identified problems:

o   In order to enable automatic join operations involving catalogues maintained in different archives, large amounts of data (potentially complete catalogues) may have to be transferred between archives using the telecommunications infrastructure.

o   To allow join operations between different catalogues these must be homogenized at least wrt. the join attributes, a task which if carried out dynamically can be very CPU demanding.

o   Furthermore it is questionable whether communication reliability is sufficient for frequent transfer of very large data sets, and finally the communication cost can not be neglected.

All of the above problems could effectively result in unacceptable response times rendering the system useless. By creating a centralized database of the archive directories (the logs and catalogues) and homogenizing a limited number of key attributes, these problems could be avoided. It does raise questions about the validity of the original fully distributed architecture, and the extent to which the pilot project confirms it.

The above problems can ideally be solved by having a dedicated communications infrastructure providing high-speed data transfer and making use of known distributed database techniques, e.g. semi-joins and semi-outer-joins, to limit the amount of data which is actually transferred between the individual archive nodes. However this will also require a very intelligent query optimizer, so at present it seems realistic to conclude that a certain directory information must be available in a few central access points, while the actual data can be maintained under the control of the archive in-

stitutions.

The ESIS Bibliographic database (ESISBIB) has been created to host bibliographic data dedicated to the two disciplines in the current scope of the QE. Rather than a traditional bibliographic DBMS, it has been implemented using a full text retrieval system, and as such serves the purpose of evaluating this approach against the more traditional DBMS'.

## 3.2   The User Platform

The User Platform provides two primary functions: query construction and result handling. These are concerned with the construction and submittal of queries and the subsequent presentation and handling of results.

### 3.2.1   Query Construction

As mentioned previously, one of the underlying assumptions of the ESIS QE system is that it is possible to provide a common data model for a given domain, covering all the involved archives and databases of that domain. The purpose of this model is to serve as the conceptual model through which the user expresses his information needs, i.e his queries. This idea has been pursued within ESIS along several parallel lines and before the QE project such models were established for the scientific domains.

These models, the Astronomy Query Language (AQL) and the Space Physics Query Language (SPQL), expressed as entity relationship diagrams, provide the vocabulary of discourse, i.e. the entities, attributes and relations used to model the two domains. The original AQL model may be found in [Albrecht 91], the SPQL may be found in [Giaretta 90], figure 2 depicts the upper most layer of the AQL model, as it looks today after further clarifications and refinements.

The user models are not normalized relational database model, but contain several elements requiring specific interpretation and query facilities. Among these are many-to-many relationships, implicit relations, and hierarchical attributes.

In relational databases a relational is implemented by having common attributes in the related tables (keys and foreign keys). Using SQL the user must know which attributes to use when expressing a join across a relation between two tables. With the implicit relations this is left to the ESIS QE system.

794

Rather than writing:

Select   Identifier, Title, Author
From    Observational_Entity, Publication
Where   Publication_Code = Asc.Number and Author
         = "Smith"

The user may write:

**Search for**
     Observational_Entity show Identifier
**Restricted By**
     Referenced_In Publication show Title, Author
**Which Has**
     Author = "Smith"

using the ESIS Conceptual Query Language (ECQL).

In this example the user need not know how the Reference_In relation is solved, and indeed it may be solved differently on different archives. Furthermore the relation may not be solvable on a single archive but only across two or more archives. In formulating a given query using the domain model the system pilot system currently provides two means of interaction:

**Query By Command,**
i.e. where the user types a query using the ECQL and the entities, attributes and relations of the discipline specific domain model. In : **Search for** *Publication* show *Title*, *Author* **Restricted By** *Written By Scientist* which has *Name* = 'Jones' show *Address*, the underlined terms denote entities, attributes and relations in the domain model.

**Query By Menu,**
the entity relationship diagram is represented in a set of menus. By selecting entities and relationships from these menus it is possible for the user to construct a valid query without having to know the precise format of the query language.

Figure 3 depicts part of the current ESIS QE user interface, where both the Query By Command and the Query By Menu windows are open.

It is furthermore planned that the ESIS QE will support:

**Query By Diagram,**
which is directly related to the graphical representation of the entity relationship diagram. From a graphical representation the user uses a mouse to select entities

and relationships from the diagram, thereby formulating a skeleton query which is refined by adding conditions and projection clauses[GESI].

**Query By Form,**
provides a set of predefined query skeletons expressed as a set of forms which is refined by adding selection clauses for the involved attributes. In terms of the datamodel the Query By Form means of interaction provides a set of views. In the pilot ESIS QE system this means of interaction will be the least emphasized, consisting mainly of a set of predefined query templates.

### 3.2.2 Special Functions

In support of the discipline specific requirements the user may apply a number of special functions. In the Astronomy domain the core ones are concerned with the correlation of star coordinates, enabling the application fuzzy joins between star or observational catalogues. Example functions are:

**In_Cone** (Coord_1, Coord_2, Radius)
This function take as input the coordinates of two observations and a radius, and returns true of Coord_2 is within the cone defined by Coord_1 and the radius. The function may be used to search in individual catalogues or to join two catalogues by coordinates.

**In_Box** and **Nearest_Object** allow similar functionality to be expressed.

In Space Physics similar functions are available to compare and join observation series ordered by time, in order to investigate time and location specific events across missions.

### 3.2.3 Result Handling

Once the user has completed a query, it is submitted for execution. This is handled by the Query Executor, which will formulate the required subqueries, issue them to the Local Query Servers, and when a result is available inform the user.

Query results will either be table oriented or consist of data and descriptor files for images, spectra, raw data and the like. These results are initially put into the Personal Database of the user, hosted at the Access Point. Using the result handling facilities, the user may select any result and decide whether and how to have it presented on his User Machine, or

decide to transfer it to e.g. the Correlation Environment. At present the system supports table presentation and very basic image presentation facilities, implemented using public domain packages. Figure 4 shows the Table Result Window.

### 3.3 The Access Point

The Access Point hosts the Query Executer and Result Management components. The primary task is to decompose user queries, submit them to the archive nodes, collect and integrate the results, and manage the results belong to various users.

The decomposition process can be described in the following steps:

1. The query is analyzed, and implicit join conditions are made explicit and many to many relations are solved using a normalized conceptual data model.

2. The query is converted into an extension query, i.e. the mapping from entities and attributes of the query is used in order to find the set of extension objects (i.e. archives, entitites and attributes) which contain information relevant to the query components.

3. The minimal set of archives necessary and sufficient to solve the query is identified.

4. The preferred set of archives for solving the query is selected.

5. The query is decomposed into sub-queries. A sub-query either address one specific archive and is expressed on the logical data model for that archive, or the sub-query combines the results from previous sub-queries. A present only very basic optimization is done in order to minimize data transfer between archives and the access point.

6. Each sub-query is sent to the appropriate archive local query server where the final mapping to the physical model and DBMS is made, and a host query is issued.

7. The result from the archive is collected by the local query servers, converted to a common format and transferred to the Access Point, where further script commands or sub-queries in the standard algorithm may combine results or operate on these results.

8. The final result is obtained, stored in the users Personal Database on the Access Point, and a notification indicating the size of the result is returned to the User Platform.

In addition to the modules directly related to the primary functionality of the ESIS QE, a major part of the software in the access point is required to manage the client server interprocess communications, the user and result management and general low level process handling.

The interprocess communication is based on DECnet and TCP/IP. On top of these a dedicated communications layerhas been implemented (the Application Communication System) which provides a common Application Programming Interface (API) accross the various platforms, i.e. SUN/UNIX and VAX/VMS, and a gateway allowing interprocess communication between a SUN/UNIX user platform running the TCP/IP protocol and VAX/VMS DECnet based environment.

### 3.4 The Archive Nodes

The physical access to an archive is performed by a dedicated Local Query Server which solves the query in the following steps:

1. The query received from the query executer is mapped on to the local schema.

2. Special functions are analyzed to see whether they are supported by the archive. Currently only functions supported by the archive is handled but with some limitations the special scientific functions could be applied in the LQS external to the archive.

3. A connection to the archive is established and a series of commands representing the query is issued.

4. The result is received and converted into the format of the common RDBMS, currently ORACLE, before it is returned to the access point.

Seen from the Query Executor each archive node is in principle a relational database with standard RDBMS capabilities, i.e. search and join functions. This implies that result combination involving temporary results can take place at any archive node. However, due to the very diverse nature of the host DBMS, ranging from commercial RDBMS, over archive

796

specific DBMS's to the ESISBIB full text DBMS, this can not be fully achieved within the current project.

## 4. Current Status and Conclusion

The coding of the ESIS QE started in January 1992, and at the current moment (August 1992) a first development release is forthcoming.

As the system becomes more stable it is the intention to release it, initially to a small set of select users, prior to the public evaluation and distribution.

The evaluation of the Pilot Project will take place in the first and second quarter of 1993, at which time the devlopment of the ESIS QE pilot project will be completed.

## 5. Acknowledgements

The authors wish to thank everyone within the ESIS team at ESRIN, within the industrial Consortium and at the various scientific institutions involved in the ESIS project. They have all in one way or another contributed to this paper and the results achieved thus far..

## 6. References

[Albrecht]
Miguel A. Albrecht
ESIS A science Information System
in "Databases & On-Line Data in Astronomy"
ed. M. Albrecht and D. Egret.
Kluwer Academic Press, 1991.

[Giaretta]
D.L. Giaretta et al
Rutherford Appleton Laboratories and
University of Sheffield
ESIS Draft Final Report The Space Physics Query Language

[GESI]
GESI
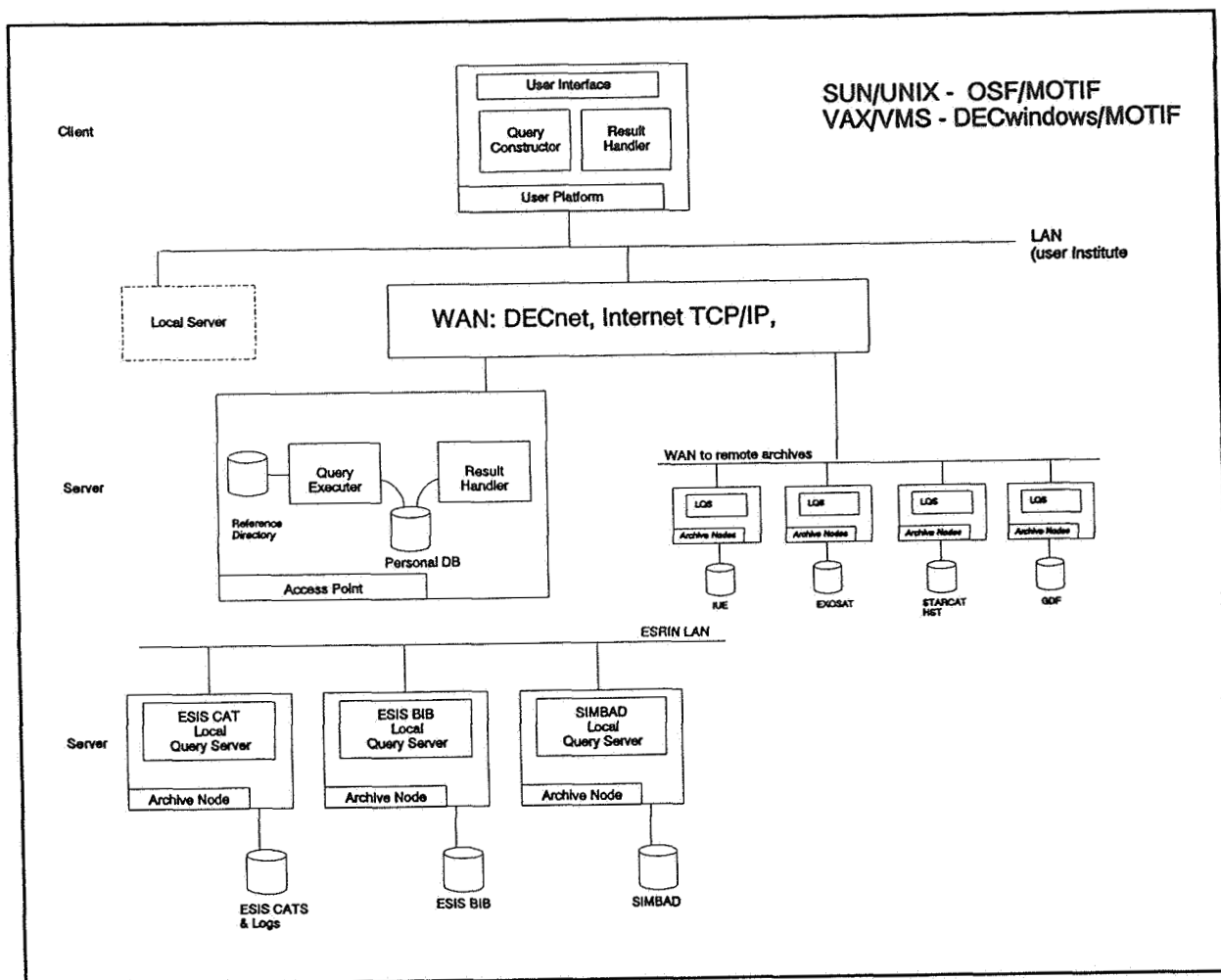The KIM Conceptual Query Language
GE-05-06-008
July 16, 1990

Client

User Interface

Query Constructor | Result Handler

User Platform

SUN/UNIX - OSF/MOTIF
VAX/VMS - DECwindows/MOTIF

LAN
(user Institute

Local Server

WAN: DECnet, Internet TCP/IP,

Server

Query Executer

Result Handler

Reference Directory

Personal DB

Access Point

WAN to remote archives

LQS | LQS | LQS | LQS

Archive Nodes | Archive Nodes | Archive Nodes | Archive Nodes

IUE | EXOSAT | STARCAT HST | GDF

ESRIN LAN

Server

ESIS CAT Local Query Server

ESIS BIB Local Query Server

SIMBAD Local Query Server

Archive Node | Archive Node | Archive Node

ESIS CATS & Logs | ESIS BIB | SIMBAD

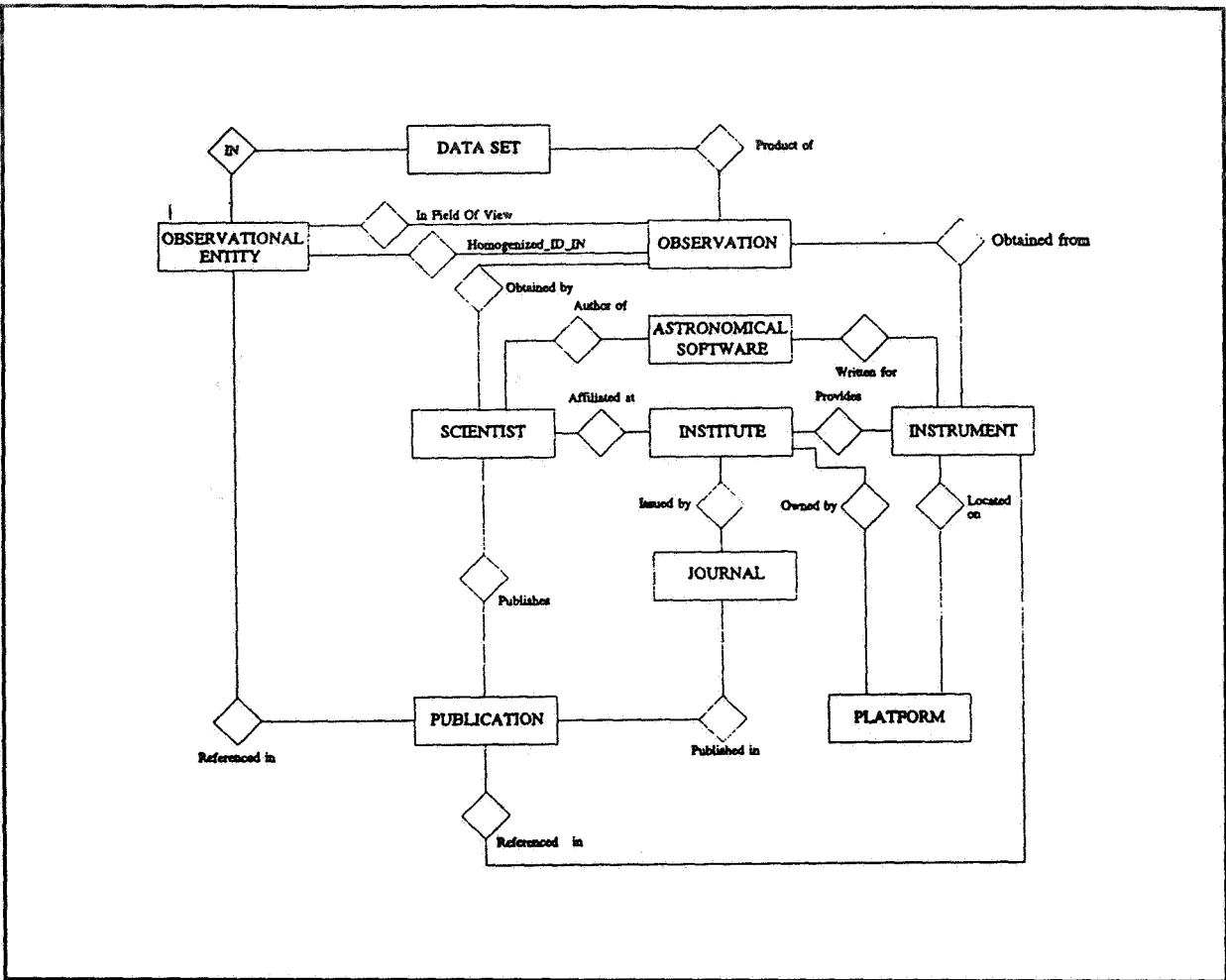Figure 1: The ESIS QE Infrastructure

798

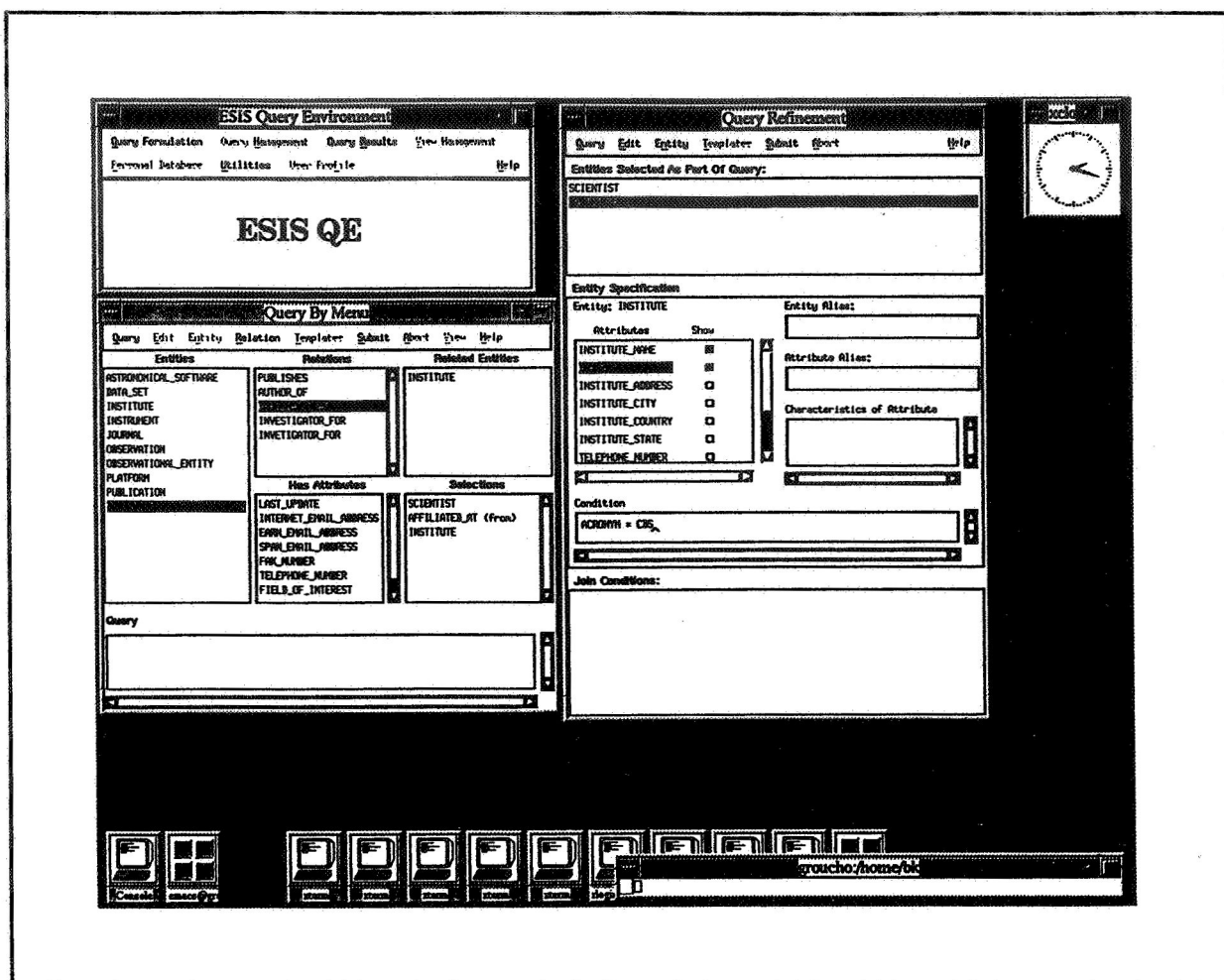Figure 2: The AQL Datamodel

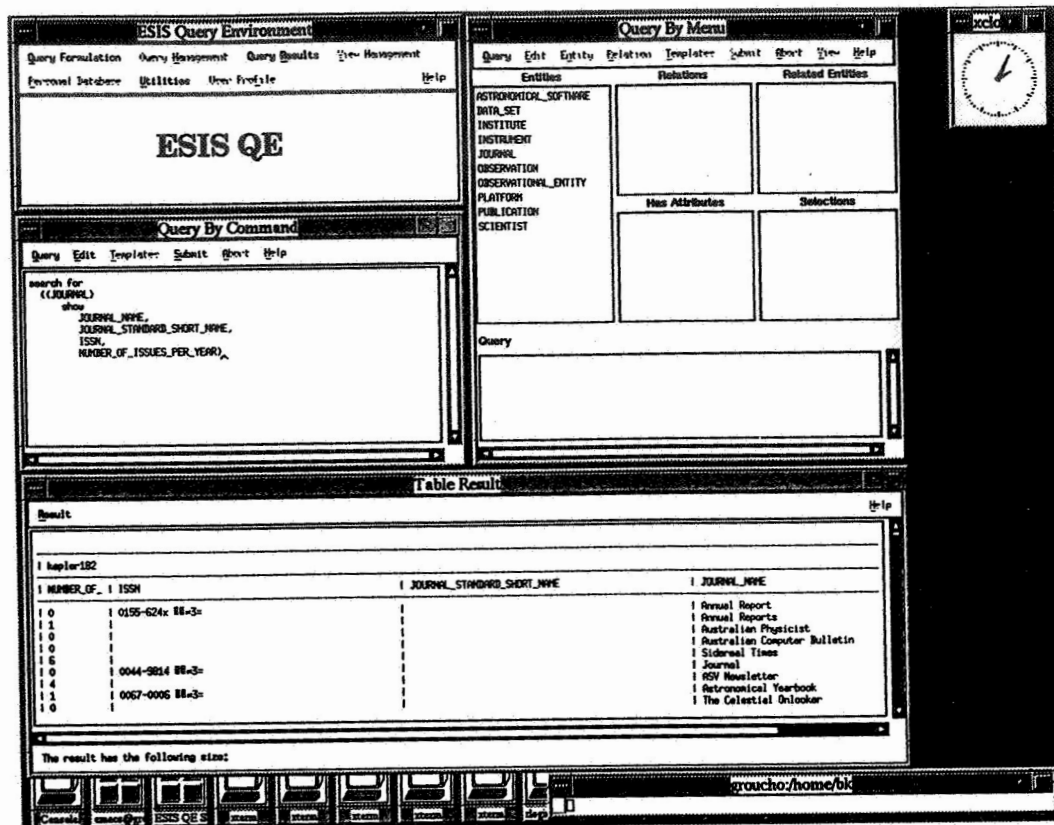Figure 3: The ESIS QE Main Menu, Query By Command and Query By Menu Windows

Figure 4: The ESIS QE Table Result Window